

# **Foundations of XML Data Manipulation**

Giorgio Ghelli

Type Systems for SSD

## Plan of the lesson

- XDuce type system (xduce-toit.2003, google:xduce)
- Tree automata (TATA, google:tree automata)
- DTD's
- XSD
- Mu-calculus and TQL Logic
- Path inclusion

## A Query

```
for $b in $doc /bib/book,  
    $a in $b /author  
where $b /@year > 2000  
return booksbyaut [$a,  
                  for $bb in $doc /bib/book  
                  where $a isin $bb /author  
                  return $bb/title  
]
```

## Its Type

- Type:

Bookbyaut[ author[String], title[String]\* ]\*

- If:

```
$doc: bib[ book [ title[String],  
                    author[String]*  
                ]*  
]
```

## Types for SSD

- XDuce type system:

T ::= B	base types
m[T]	tree types
T, T	forest concatenation types
0	empty tree singleton type
T   T	union type
X	m[ ]-guarded recursion
T*	Kleene star type

## Recursion

- Section = para[] | section[Section\*]
- Illegal:
  - paraList = para[] | (para[], paraList)
- Equivalent to:
  - paralist = para[], (para[])<sup>\*</sup>
- Illegal:
  - X = a[], X , b[] | 0

## Type rules

$$\frac{}{0 : 0} \quad \frac{t : T}{m[t] : m[T]} \quad \frac{t : T \quad u : U}{t, u : T, U}$$

$$\frac{t : T}{t : T|U} \quad \frac{}{0 : T^*} \quad \frac{t : T \quad u : T^*}{t, u : T^*}$$

$$\frac{X=T \vdash t : T}{X=T \vdash t : X}$$

## Not so different?

- Is  $T, T'$  the same as  $T \times T'$ ?
- Is  $T^*$  the same as  $\text{List}(T)$ ?
- Is  $a[T], b[U], c[V]$  the same as  $\text{record}[a:T, b:U, c:V]$ ?

## $T, T'$ vs. $T \times T'$

- $(T, T'), T'' = T$ ,  $(T', T'') = T$ ,  $T', T''$
- $T, 0 = T$  but  $T \times 1 \neq T$
- $\langle t, u \rangle : T \times U$  iff  $t : T$  and  $u : U$
- It may be that  $t, u : T, U$  but not  $t : T$  and  $u : U$ :
  - $(a[], b[], c[]): (a[], (b[], c[]))$
  - $a[], (b[], c[]): (a[] | a[], b[]) , (b[], b[] | c[])$
  - Type checking similar to regular-language testing
- $T \times U < T' \times U'$  iff  $T < T'$  and  $U < U'$
- It may be that  $T, U < T', U'$  but  $T \not\leq T'$  and  $U \not\leq U'$
- Subtyping defined as language inclusion, and checked as automata inclusion

## The semantic intuition

- At run time, a value of type  $T \times U$  is represented as `pair(t,u)`
- A value of type  $T \times U$  is “bigger” than a value of type  $T$
- If  $t:T$  then  $t$  has NOT type  $T \times U$
- A value of type  $T, U$  has no `tag(.,.)` structure
- If  $t:T$  and  $0:U$  then  $t: (T,U)$

`a[T],b[U]` vs. `tup[a:T] + tup[b:U]`

- Record concatenation:
  - $\text{tup}[a:T] + \text{tup}[b:U] = \text{tup}[a:T, b:U]$
- Usually  $T+U$  only defined when  $T$  is a simple-record type  $\text{tup}[a:T]$ ;  $T, U$  accepts any type as  $T$ , including  $(a[] \mid 0)$

## $T^*$ vs. $\text{List}(T)$

- $T^* = \mu TS. 0 \mid (T, TS)$
- $\text{List}(T) = \mu LT. 1 \mid (T \times LT)$
- Hence:
  - $T^* = T^{**}$
  - $\text{List}(T) \neq \text{List}(\text{List}(T))$
  - $T < T^*$ : no List tag at run-time
  - $T \not\in \text{List}(T)$

## Guarded recursion

- Regular grammars:
  - $X = a_1 \mid a_2 \mid a_3.X_1 \mid a_4.X_n$
- XDuce “linear” recursive types:
  - $X = a_1[] \mid a_2[] \mid a_3[X_1] \mid a_4[X_2]$
  - Correspond to automata
- XDuce tree-like recursive types
  - $X = a_1[] \mid a_2[] \mid a_3[X_1, X_2] \mid a_4[X_3, X_4]$
  - Correspond to tree automata

## Horizontal and Vertical RegExps

- Horizontal:
  - $X = a[](b[])^* \mid (a[])^*b[]$
- Vertical
  - $X = a[Y] \mid W$
  - $Y = 0 \mid b[Y]$
  - $W = b[] \mid a[W]$
- Tree automata because of Vertical
- *Regular* tree automata because of Horizontal

## Tree Automata

- $(A, Q, R, F)$  with
  - $F \subseteq \text{Lists}(Q)$
  - $R \subseteq A \times \text{Lists}(Q) \times Q$  (set of  $a[q_1, \dots, q_n] \rightarrow q$  rules)
- A run:
  - Substitute  $a[]$  with  $q$  ( $a^q[]$ ) if  $a[] \rightarrow q \in R$
  - Substitute  $a[q_1, \dots, q_n]$  with  $q$  if  $a[q_1, \dots, q_n] \rightarrow q \in R$
  - Accept  $t_1, \dots, t_n$  if rewritten as  $q_1, \dots, q_n \in F$

## Unranked Tree Automata

- Ranked Tree Automata:
  - Rules are like  $a_2[q_1, q_2] \rightarrow q_3$ : for each binary symbol, I need  $2^{3*|Q|}$  rules at most
- But:  $A \times \text{Lists}(Q) \times Q$  is not finite:
  - Rules like  $a[q_1, \dots, q_1] \rightarrow q_2$
- Problem:
  - Representing  $R$  and deciding  
 $a[q_1, \dots, q_n] \rightarrow q \in R$

## Regular Tree Automata

- Regular Tree Automata:
  - For each  $a, q$ , the language (in  $Q^*$ )  $\{q_1, \dots, q_n \mid a[q_1, \dots, q_n] \rightarrow q \in R\}$  is regular
- $R$  can be represented as a function of type  
 $A \times Q \rightarrow \text{RegExp}(Q)$
- Tree automata correspond to vertical recursion
- *Regular* in RTA corresponds to horizontal regular recursion

## Unranked binary expressions

- $\text{Or}(\text{F}, \text{And}(\text{T}, \text{T}, \text{F}), \text{And}(\text{F}, \text{T}))$   
where  $\text{T} = \text{And}()$ ,  $\text{F} = \text{Or}()$
- $A = \{\text{And}, \text{Or}\}; Q = \{t, f\}$
- $R:$ 
  - $\text{Or}(f^*) \rightarrow f$
  - $\text{Or}(\ (t|f)^*, t, (t|f)^* ) \rightarrow t$
  - $\text{And}(t^*) \rightarrow t$
  - $\text{And} (\ (t|f)^*, f, (t|f)^* ) \rightarrow f$

## We like automata

- Recognize trees in linear time
- Closed by union, intersection,  
*complement*
- Emptyness is decidable
- $A <: A'$  iff  
 $A \setminus A' = A \cap \text{Co}(A')$  is empty

## DTDs

- Canonical way to describe the structure of an XML document
- Example:

```
<!ELEMENT people_list (person*)>
<!ELEMENT person (name, birthdate?, children?)>
<!ELEMENT children (person+)>
<!ELEMENT name (#PCDATA)>
<!ELEMENT birthdate (#PCDATA)>
```

## An example

- DTD:  

```
<!ELEMENT people_list (person*)>
<!ELEMENT person (name, birthdate?, children?)>
<!ELEMENT children (person+)>
```
- Document:
  - ```
<!DOCTYPE people_list SYSTEM "example.dtd">
<people_list>
  <person>
    <name>Fred Bloggs</name>
    <birthdate>...</birthdate>
    <children>
      <person><name>Jim
      </name></person>
    </children>
  </person>
  <person><name>Luis Gutierrez</name>
  </person>
</people_list>
```

## XDuce vs. DTD

- XDuce: a set of mutual recursive defs:
  - $A = a[B1^*, B2]$
  - $B1 = b[X]$
  - $B2 = b[Y]$
  - ...
- DTD: the type is identified by the label:
  - $a = a[b^*, c]$
  - $b = b[X]$
  - $c = c[Y]$

## DTD into automata

- DTD:

```
<!ELEMENT people_list (person*)>
<!ELEMENT person (name, birthdate?, children?)>
<!ELEMENT children (person+)>
```
- Automaton:
  - $F = PL$
  - $R:$ 
    - $\text{people\_list}[P^*] \rightarrow PL$
    - $\text{person}[N, BD?, C?] \rightarrow P$
    - $\text{children}[P+] \rightarrow C$
    - $\text{name}[ ] \rightarrow N \dots$

## XDuce into automata

- XDuce:
  - Paper where
  - Paper = title[], section[abstract[]], Content
  - Content = ( paragraph[] | section[Content] ) \*
- Automata:
  - $F = T, SA, (P|SC)^*$
  - R:
    - title[]  $\rightarrow$  T, abstract[]  $\rightarrow$  A,
    - section[A]  $\rightarrow$  SA
    - paragraph[]  $\rightarrow$  P
    - section[(P|SC)\*]  $\rightarrow$  SC

## XSD

- Local element types are not label-identified, but global element types are:
  - $a = a[b[X]]$
  - $b = b[a[Z]]$
  - $c = c[a[W]]$
- Not every regexp is OK (Unique Particle Attribution):
  - $a = a[b[X]^*, b[X]]$ : illegal
- In one element, label identifies type (Element Declarations Consistent):
  - $a = a[b[X], b[Y]]$ : illegal

## XSD

- Names are qualified with respect to namespaces
- XSD can specify key and keyref constraints
- Two limited forms of subtyping by name: derivation and substitution groups

## XSD Syntax

- Global element declarations:
  - element El of T
  - element El of (type [T] of {...})
  - May be local, in which case El is not a key
- Complex type definitions:
  - type T of
  - Either anonymous, or T is a key (even if local)

## XSD Assessmnet

- Assessment:
  - local validation, schema-validity assessment and infoset augmentation: Infoset -> PSVI
- PSVI contains:
  - Normalized and default values for attributes and elements
  - Type definitions for attributes and elements
  - Validation outcome

## XSD and Subtyping

- Derivation:
  - Every complex type either extends or restricts another type, starting from xsi:anyType
  - Explicit cast: the derived type can be used for validation only if a corresponding xsi:type attribute is present in the element to validate, or if the element name is in the substitution group of the expected name
- Substitution:
  - An element name may be head of a substitution group, and the other names from the group are valid where the head is required
  - The types of the group elements must be derived from the type of the head

## Semantic subtyping

- Rule-based subtyping:
  - Subsumption:  $T <: T' \Rightarrow \forall t. t:T \Rightarrow t:T'$
- Semantic subtyping:
  - Definition:  $(\forall t. t \in [[T]] \Rightarrow t \in [[T']]) \Rightarrow T <: T'$
- For example:
  - $(\text{Forall } X. X \rightarrow X) <: ? \text{ Int} \rightarrow \text{Int}$
- XSD: rule-based subtyping
- XDuce, CDuce: semantic subtyping

## DTD and $\mu$ -calculus

- [everywhere]  $A = \forall \xi (A \wedge [\downarrow]\xi \wedge [\rightarrow]\xi)$
- We extend  $\mu$  with equations:
  - $A$  where  $\$x_1 = A_1, \dots, \$x_n = A_n$
  - $A(\xi)$  where  $\xi = A_1$ 
    - is the same as  
 $A(\mu\xi. A_1)$
- Still checkable in  $O(2^n)$

## DTD and $\mu$ -calculus

- DTD:

```
<!ELEMENT people_list (person*)>
<!ELEMENT person (name, birthdate?, children?)>
<!ELEMENT children (person+)>
```

- $\mu$  with equations:

[everywhere] ( $\text{people\_list} \wedge [\downarrow] \$\text{PersonPlus}$ )  
     $\vee (\text{person} \wedge <\!\!\downarrow\!\!> \$\text{NBC})$   
     $\vee (\text{children} \wedge <\!\!\downarrow\!\!> \$\text{PersonPlus})$   
     $\vee (\text{name} \wedge [\downarrow]\text{False}) \vee \dots$

where  $\$PersonPlus = \text{person} \wedge [\rightarrow] \$PersonPlus$

$\$NBC = \text{name} \wedge [\rightarrow] \$BC$

$\$BC = (\text{birthdate} \wedge [\rightarrow] \$C) \vee \text{children}$

## TQL Logic

- Ordered TQL logic:

|               |                               |
|---------------|-------------------------------|
| $A ::= B$     | base values                   |
| $\eta[A]$     | trees ( $\eta: x$ or $n$ )    |
| $A, A$        | sequence                      |
| $0$           | empty tree singleton sentence |
| $T \vee T$    | disjunction                   |
| $\neg A$      | negation                      |
| $\mu\xi. A$   | (positive) recursion          |
| $\xi$         | recursion variable            |
| $\exists x.A$ | label quantification          |
| $\exists X.A$ | forest quantification         |
| $X$           | forest variable               |

## The actual TQL logic

- TQL data model in unordered:
  - $0|t = t$  ;  $(t|t')|t'' = t|(t'|t'')$  ;  $t | t' = t' | t$
- In TQL ordered logic:
  - $t \models (A, B)$   
iff  $\exists t', t''. t', t'' = t$  and  $t' \models A$  and  $t'' \models B$
- In TQL logic:
  - $t \models A | B$   
iff  $\exists t', t''. t' | t'' = t$  and  $t' \models A$  and  $t'' \models B$

## TQL Logic

- $F \models \text{True}$ : always ( $\text{True} = 0 \vee \neg 0$ )
- $F \models 0$  iff  $F=0$
- $F \models A | B$  iff  $\exists F', F''. F = F'|F'', F' \models A, F'' \models B$   
 $F \models m[A]$  iff  $F=m[F'], F' \models A$
- E.g.:
  - $a[0] | b[0] \models b[0] | \text{True}$  ?
  - $b[0] \models b[0] | \text{True}$  ?
  - $a[0] | b[0] \models b[0]$  ?
  - $a[b[0]] \models a[\text{True}] | \text{True}$  ?

## Other operators

- $F \models A \wedge B$  iff  $F \models A$  and  $F \models B$
- $F \models \neg A$  iff not ( $F \models A$ )
- Derived operators:
  - $A \vee B =_{\text{def}} \neg(\neg A \wedge \neg B)$
  - $A \parallel B =_{\text{def}} \neg(\neg A \mid \neg B)$
  - $m[\Rightarrow A] =_{\text{def}} \neg m[\neg A]$
- $F \models (a[\text{True}] \vee b[\text{True}]) \mid \text{True}$
- $F \models (a[0] \mid \text{True}) \wedge \neg(a[0] \mid a[0] \mid \text{True})$
- $F \models \text{author}[\Rightarrow \text{Hull}] \parallel \text{False}$

## More than types?

- Complement  $\neg A$  dualizes every other operator ( $\exists \rightarrow \forall$ ,  $\vee \rightarrow \wedge$ ,  $\mu \rightarrow \nu, \dots$ )
- Horizontal recursion is more than  $T^*$
- Quantification expresses correlation:
  - $\exists x. x[\text{True}] \mid x[\text{True}] \mid \text{True}$
  - $\exists x. .x[A] \wedge .x[A] \quad (.x[A] = x[A] \mid \text{True})$
  - $\exists x. .x[A] \mid .x[A] \quad (.x[A] = x[A] \mid \text{True})$
- Logic can express key constraints:
  - $\neg \exists X. .\text{book}[\cdot.t[X]] \mid .\text{book}[\cdot.t[X]]$

## Decidability

- Quantification makes emptiness undecidable
- Quantification makes model-checking (type-checking) PSpace-complete
- Model-checking is often doable in practice

## Path containment

## Path containment

- As binary relation:  $\text{sub}_2$ 
  - $p \text{ sub}_2 q \Leftrightarrow (m \text{ } p \text{ } n \Rightarrow m \text{ } q \text{ } n) \Leftrightarrow [[p]] \subseteq [[q]]$
- Starting from the root:  $\text{sub}_1$ 
  - $p \text{ sub}_1 q \Leftrightarrow (\text{root } p \text{ } n \Rightarrow \text{root } q \text{ } n)$
- Boolean containment, starting from the root:
  - $t \models p$ : matching  $p$  against the root of  $t$  yields non-empty result
  - $p \text{ sub}_0 q$  iff  $t \models p \Rightarrow t \models q$

## Notions of containment

- If we restrict to child/desc,  $\text{sub}_2$  e  $\text{sub}_1$  are equivalent
- In the presence of predicates,  $\text{sub}_1$  can be mapped to  $\text{sub}_0$  :
  - $p \text{ sub}_1 q \Leftrightarrow p[x] \text{ sub}_0 q[x]$ , where:
    - $p[x]$  adds a  $\text{child}::x$  condition to the selection node, and  $x$  is fresh (Miklau-Suciu PODS 02)

## Complexity for PositiveXPath

- PTime:
  - No disjunction, 2 of  $// []^*$  but not all:  
 $\text{XP}(/, //, *)$ ,  $\text{XP}(/, [], *)$ ,  $\text{XP}(/, //, [])$ ,  $\text{XP}(/, //) + \text{DTD}$ ,
- coNP:
  - $\text{XP}(/, |)$ ;  $\text{XP}(//, |)$ ;
  - $\text{XP}(/, []) + \text{DTD}$ ,  $\text{XP}(//, []) + \text{DTD}$
  - $\text{XP}(/, //, [], *)$ ;
  - $\text{XP}(/, //, [], *, |)$  (becomes PSPACE if the alphabet is finite);
- ExpTime:
  - $\text{XP}(/, //, |) + \text{DTD}$
  - $\text{XP}(/, //, [], |, *) + \text{DTD}$

## Path inclusion and $\mu$ -calculus

- Let  $<<\text{p}>>$  be the translation of a path and  $<<\text{s}>>$  the translation of a schema:
  - $E, L, m \models <<\text{s}>>$  if  $E, L$  satisfies  $s$
- $p \text{ sub}_2 q$ :
  - Valid ( $<<\text{p}>> \Rightarrow <<\text{q}>>$ )
  - i.e., for any  $E, L, m, n$ :  
 $E, L, i \rightarrow n, m \models <<\text{p}>> \Rightarrow <<\text{q}>>$
- $p \text{ sub}_2 q$  under  $s$ :
  - For any  $E, L, m, n$ :  
 $E, L, i \rightarrow n, m \models <<\text{p}>> \wedge <<\text{s}>> \Rightarrow <<\text{q}>>$
- $<<\_>>$  is linear  $\Rightarrow$  inclusion is  $O(2^n)$  for NavXPath